# Mitigating Cross-Database Differences for Learning Unified HRTF Representation

*Yutong Wen*, *You Zhang, Zhiyao Duan*

University of Rochester, NY, USA

ywen6@u.rochester.edu

IEEE WASPAA 2023

# Head-related transfer functions (HRTFs)

HRTFs encode **spectral changes** of sound from source to listener's ears.

HRTFs are fundamental to virtual auditory displays.
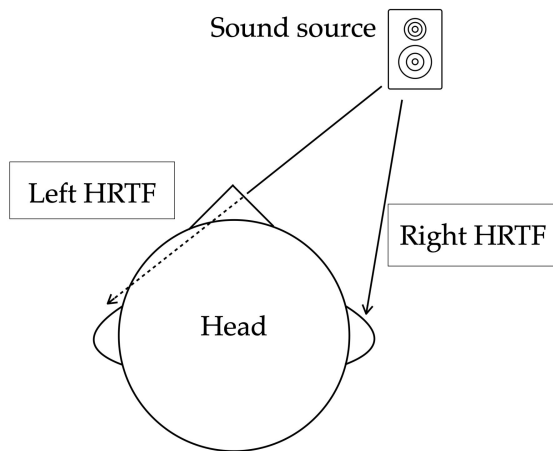
It is **hard to measure** personalized HRTFs.



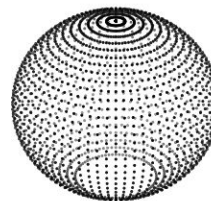Figure from: https://ieeexplore.ieee.org/document/7099223

# Existing databases

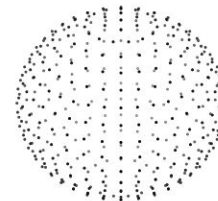| Database Information | ARI | ITA | Listen | Crossmod | SADIE II | BiLi | HUTUBS | CIPIC | 3D3A | RIEC |
|---|---|---|---|---|---|---|---|---|---|---|
| # Subjects | 97 | 48 | 50 | 24 | 18 | 52 | 96 | 45 | 38 | 105 |
| # Positions | 1550 | 2304* | 187 | 651 | 2818* | 1680 | 440 | 1250 | 648 | 865 |
| Source Distance (m) | 1.2 | 1.2 | 1.95 | 1.0 | 1.2 | 2.06 | 1.47 | 1.0 | 0.76 | 1.5 |

Small number of subjects in each database

- Hard to model such high-dimensional data

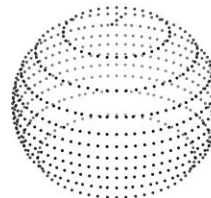Different spatial sampling schemes

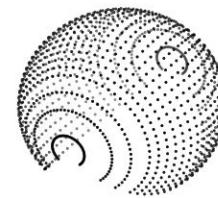- Hard to do mix-database training
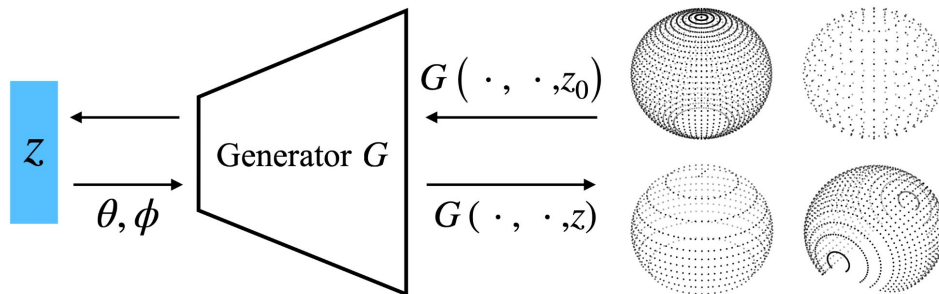
ITA    HUTUBS

Crossmod    CIPIC

# HRTF field for personalized HRTF modeling

We proposed ***HRTF field*** to alleviate the differences of spatial sampling schemes, enabling mix-database training [1].

The basic idea is to view an HRTF as a discrete sample from the underlying continuous-space HRTF and **model the continuous function** directly.
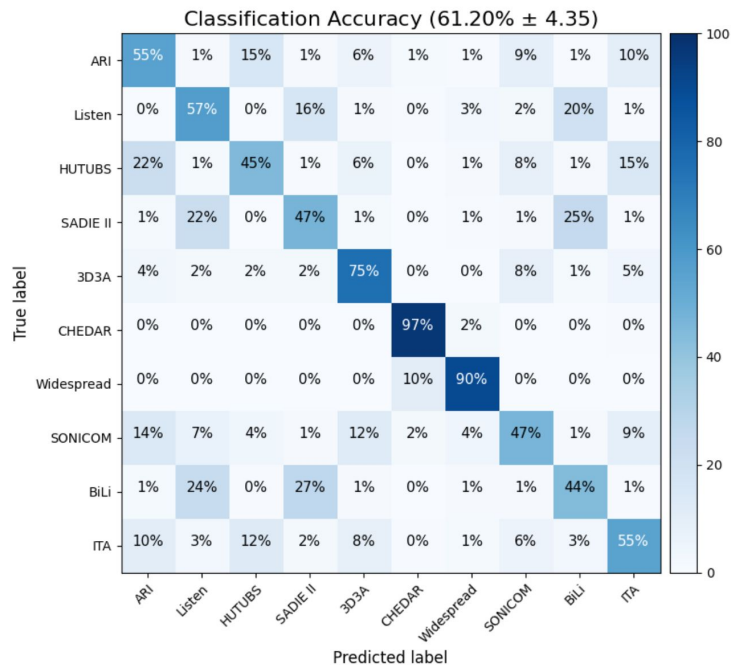


[1] Zhang, You, Yuxiang Wang, and Zhiyao Duan. "HRTF field: Unifying measured HRTF magnitude representation with neural fields." *ICASSP 2023*.

4

# Differences beyond spatial sampling schemes

A recent study shows that there are other significant differences across HRTF databases.

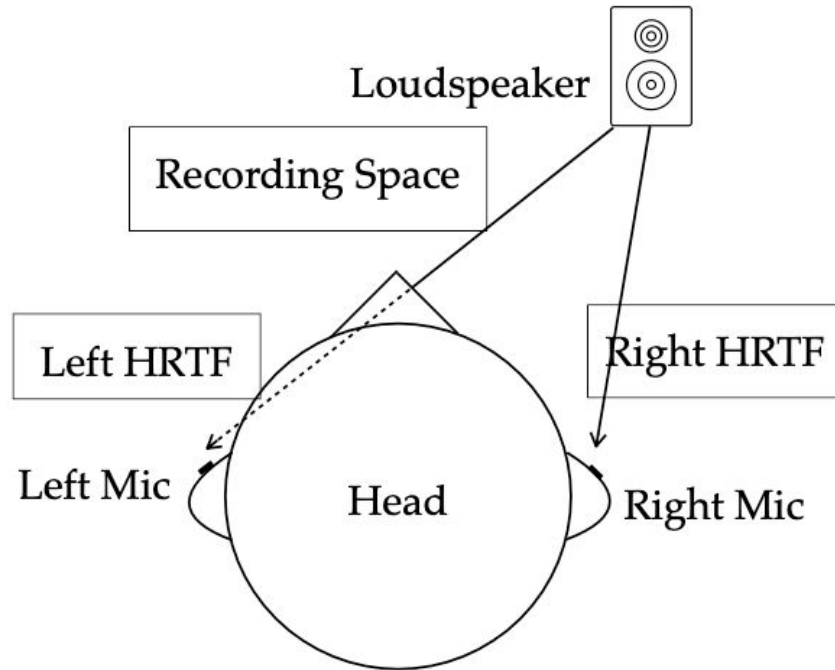This would hinder the training process.



Pauwels, Johan, and Lorenzo Picinali. "On the relevance of the differences between HRTF measurement setups for machine learning." *ICASSP 2023*.
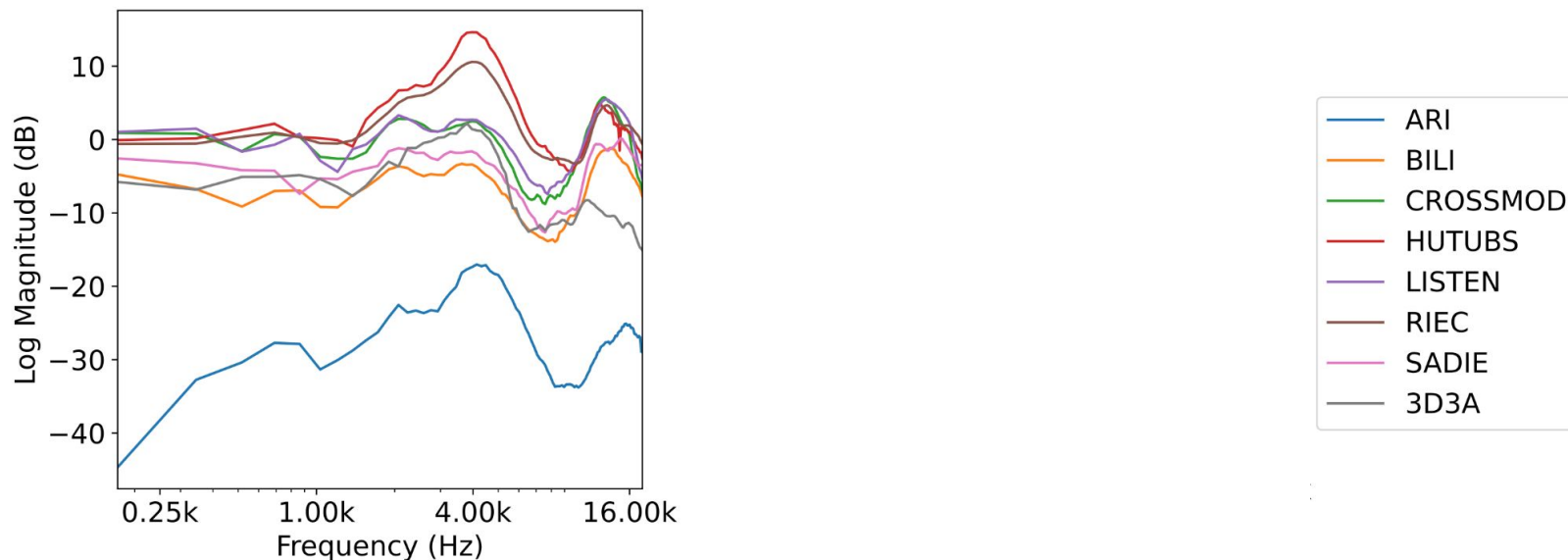
# Investigating cross-database differences

Frequency response of:

1. Different loudspeakers,

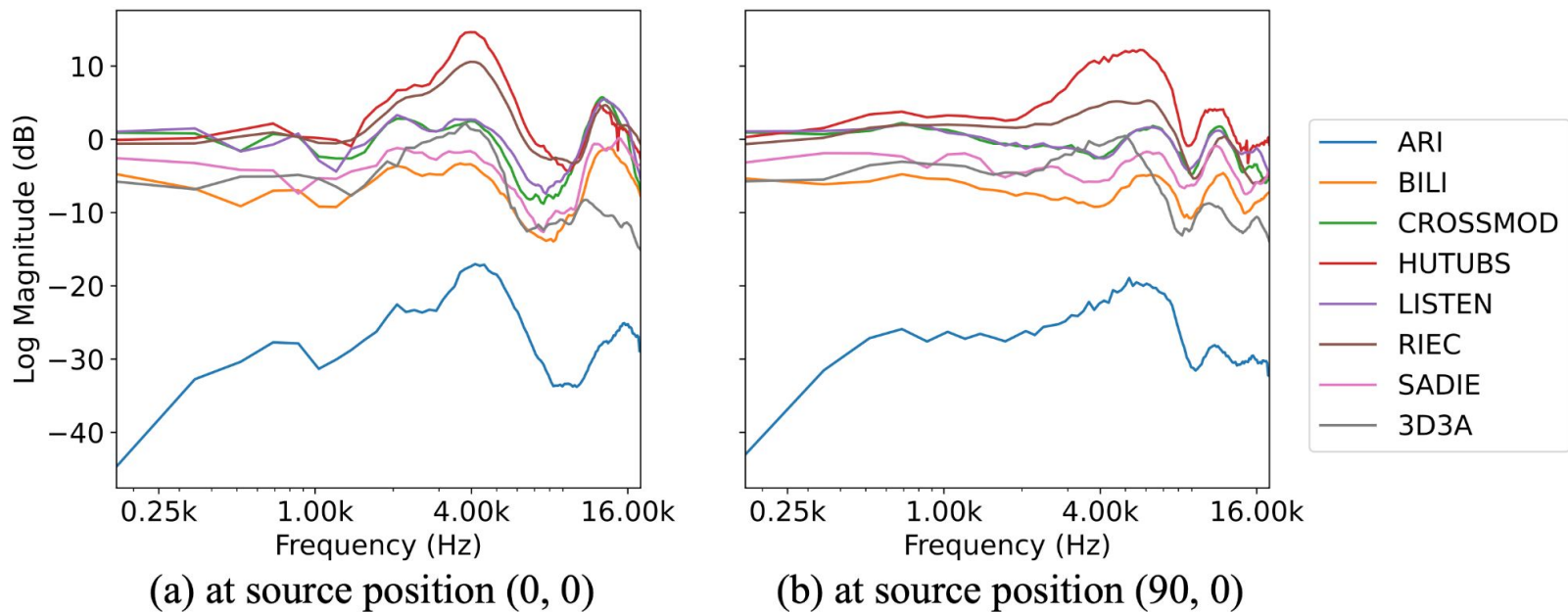2. Recording spaces,

3. Microphones.

# Average HRTFs across subjects



(a) at source position (0, 0)

There are systematic differences in measurement system responses at each source position.

# Average HRTFs across subjects



(a) at source position (0, 0)   (b) at source position (90, 0)

These position-dependent systematic differences in measurement system responses need to be removed.

# Our method to normalize HRTFs

Unnormalized HRTF magnitude

$$HRTF_{\text{normalized}}(\theta, \phi) = \frac{Y(\theta, \phi)}{HRTF_{\text{avg}}(\theta, \phi)}$$
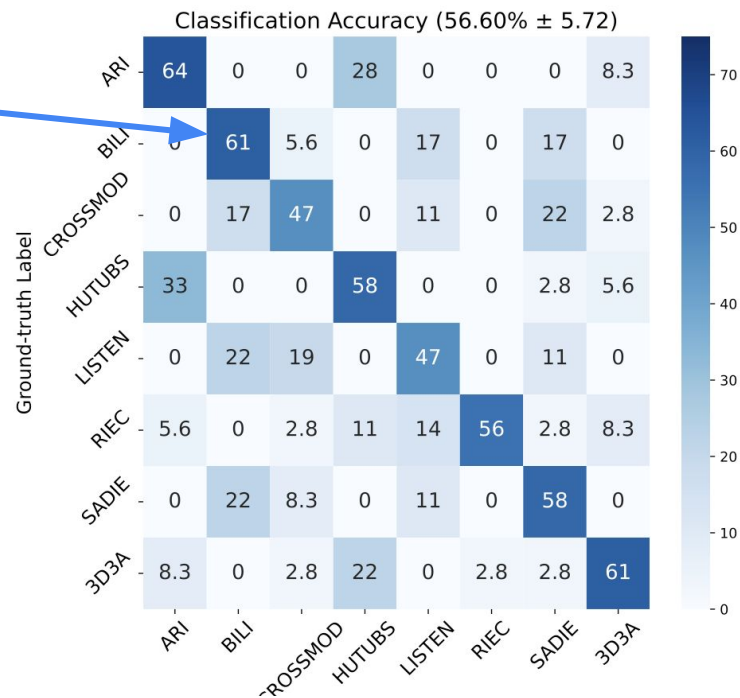
azimuth

elevation

Average HRTF magnitude across subjects

# HRTF database classification before normalization

Experimental setup:

- **Frequency range**
  - 200Hz to 18kHz
  - 104 frequency bins
- **Total 144 subjects**
  - 18 (the smallest size dataset) times 8
  - 432 HRTFs = 18 (subjects)
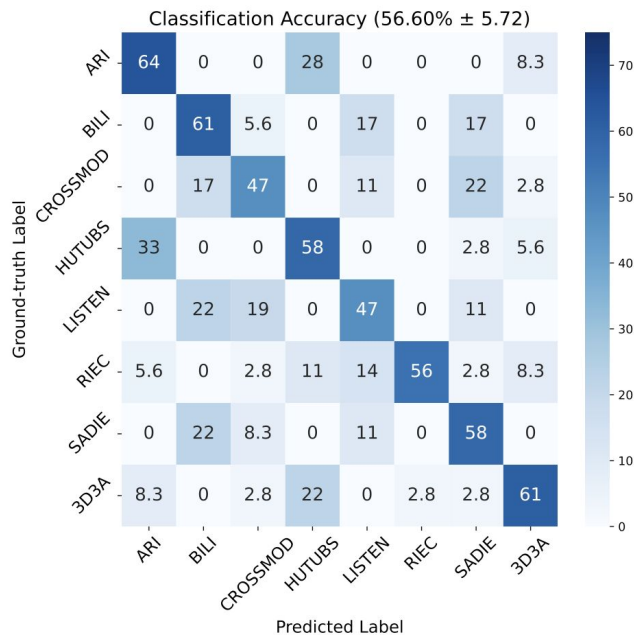    * 12 (common positions) * 2 (ears)
- **Model: kernel SVM**

Numbers are percentage
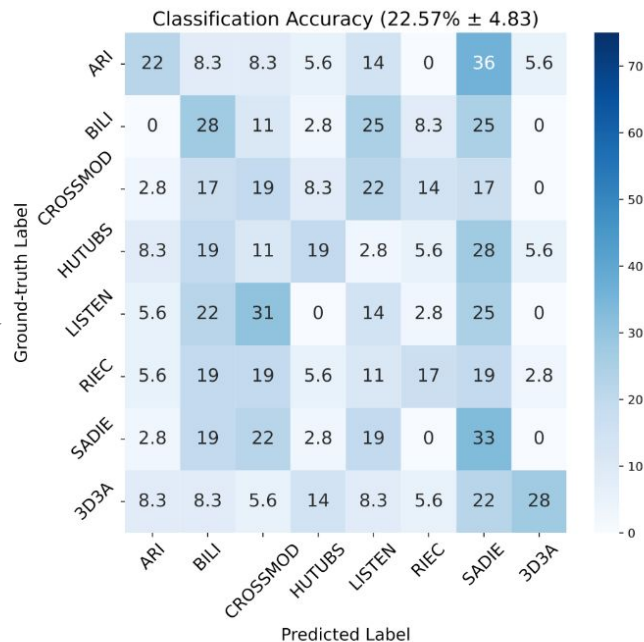


Classification Accuracy (56.60% ± 5.72)

SVM can easily tell where the HRTFs originate from.

Pauwels, Johan, and Lorenzo Picinali. "On the relevance of the differences between HRTF measurement setups for machine learning." *ICASSP 2023*.

# Our normalization successfully confuses SVM classifier



$$HRTF_{\text{normalized}}(\theta, \phi) = \frac{Y(\theta, \phi)}{HRTF_{\text{avg}}(\theta, \phi)},$$
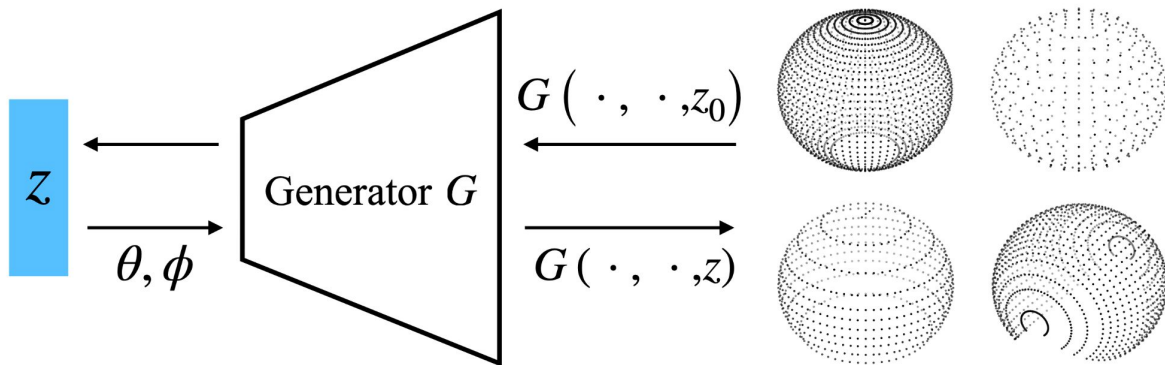
# Experiment with *HRTF field*

Use *HRTF field* to do cross-database reconstruction.

We train the generator with multiple databases combined.

Inference HRTFs from another database.



$G(\ \cdot\ ,\ \cdot\ ,z_0)$

$z$

Generator $G$

$\theta, \phi$

$G(\ \cdot\ ,\ \cdot\ ,z)$

Zhang, You, Yuxiang Wang, and Zhiyao Duan. "HRTF field: Unifying measured HRTF magnitude representation with neural fields." *ICASSP 2023*.

# Evaluation on log-spectral distortion (LSD)

Ground-truth

Predicted

# frequency bins

# spatial locations

Frequency index

$$LSD(H, H') = \sqrt{\frac{1}{PN} \sum_{\theta, \phi} \sum_{n} \left(20 \log_{10} \left| \frac{H(\theta, \phi, n)}{H'(\theta, \phi, n)} \right| \right)^2}$$

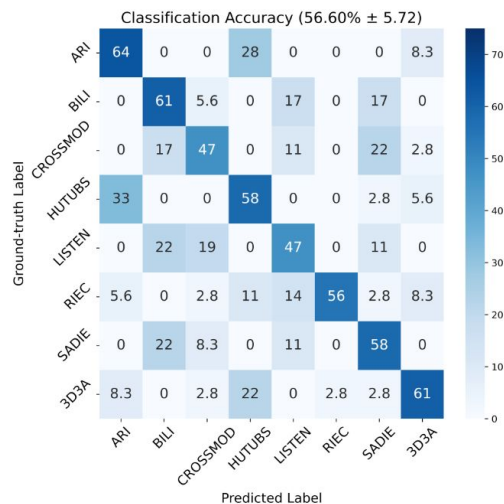# Our normalization improves cross-database reconstruction

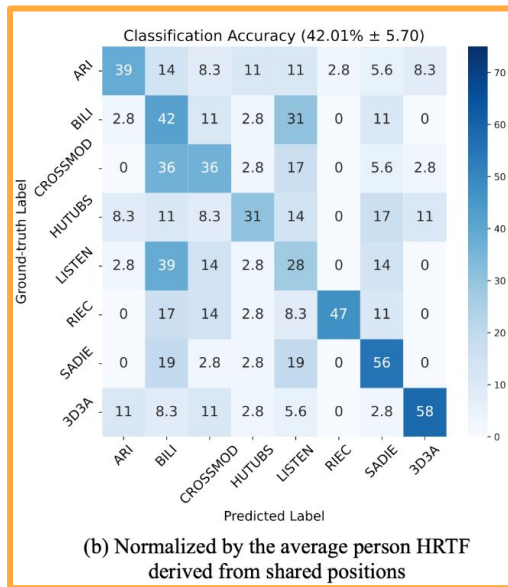Training sets are denoted with △

Testing sets are denoted with ◯

HRTF field with our normalization method

| Experiments | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| ARI | ◯ | △ | | △ | △ |
| ITA | | | | △ | △ |
| Listen | △ | | ◯ | △ | △ |
| Crossmod | △ | △ | △ | △ | △ |
| SADIE II | △ | | △ | △ | △ |
| BiLi | △ | △ | △ | △ | △ |
| HUTUBS | | △ | | △ | ◯ |
| CIPIC | | | | △ | △ |
| 3D3A | | | | △ | △ |
| RIEC | | ◯ | | ◯ | △ |
| HRTF field [15] | 7.47 | 5.54 | 4.31 | 4.43 | 5.01 |
| **Our proposed** | **4.69** | **4.82** | **3.89** | **3.73** | **4.04** |

Yutong (Cooper) Wen       Mitigating Cross-database Differences for Learning Unified HRTF Representation       WASPAA 2023
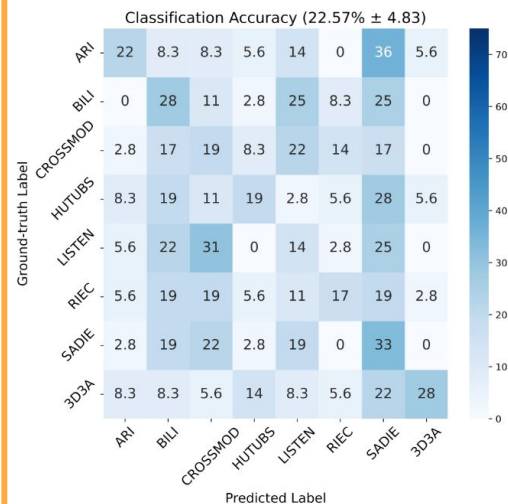
# Ablation study



(a) Unnormalized data

(b) Normalized by the average person HRTF derived from shared positions

(c) Normalized by the average person HRTF derived from individual position

| | | | | | |
|---|---|---|---|---|---|
| HRTF field [15] | 7.47 | 5.54 | 4.31 | 4.43 | 5.01 |
| **Our proposed** | **4.69** | **4.82** | **3.89** | **3.73** | **4.04** |
| w/o position dependency | 5.61 | 5.32 | 4.32 | 4.00 | 4.89 |
| w/o ear dependency | 5.11 | 5.11 | 3.98 | 3.94 | 4.67 |

# Takeaways

- There are **position-dependent** systematic differences across HRTF databases.

- It is effective to normalize these differences using **average person HRTFs from individual positions.**

- Our proposed normalization method is promising to benefit many machine learning based methods for HRTF research.

ywen6@u.rochester.edu

Full text

Code